

Zpravodaj Československého sdružení uživatelů TeXu

Janka Chlebíková

Ako rozdělit (slovo) Československo

Zpravodaj Československého sdružení uživatelů TeXu, Vol. 1 (1991), No. 4, 10–13

Persistent URL: <http://dml.cz/dmlcz/148815>

Terms of use:

© Československé sdružení uživatelů TeXu, 1991

Institute of Mathematics of the Czech Academy of Sciences provides access to digitized documents strictly for personal use. Each copy of any part of this document must contain these *Terms of use*.



This document has been digitized, optimized for electronic delivery and stamped with digital signature within the project *DML-CZ: The Czech Digital Mathematics Library* <http://dml.cz>

Ako rozdeliť (slovo) Československo.

Každý, kto používal \TeX na písanie v slovenčine či češtine v dobe, kedy \TeX ešte nebol prispôsobený na naše československé (česko-slovenské, české a slovenské)² pomery, musel sa vysporiadať s dvoma nákladnými problémami:

1. Akcentované znaky – špeciálne ĺ, ě, ť, ($\backslash v$ totiž láva ť, s čím sa nemožno uspokojiť). Nemožno tiež zabudnúť na „naše“ úvodzovky, ktoré sa líšia od „anglických“.
2. Rozdeľovanie slov – neprispôsobené pôvodné verzie \TeX u rozdeľujú slová podľa pravidiel anglickej gramatiky, vôbec nerozdeľujú akcentované slová a ani ich nedovoľujú pridávať do slovníka výnimiek $\backslash\text{hyphenation}\{\dots\}$.

Prvý problém je možno v súčasnosti veľmi pohodlne riešiť (odhliadnuc od pamäťových nárokov) napr. použitím 8-bitových fontov od P. Nováka. Pred dokončením sú už i 8-bitové DC-fonty, obsahujúce väčšinu znakov z najrozšírenejších európskych jazykov (i naše ě, ě, ĺ, f). Tieto fonty by sa mali stať medzinárodným štandardom. Navyše všeobecne rozšírený $\text{em}\text{\TeX}$ poskytuje pohodlnú možnosť konverzie akcentovaných znakov do ľubovoľných kódov, takže editovanie si môžeme spríjemniť napr. použitím známeho programu od Kamenických. Podľa mojich osobných skúseností ani posledná verzia PC \TeX u neumožňuje tak jednoduchú konverziu ako $\text{em}\text{\TeX}$ (Američanov až natalko netrápia naše akcentované znaky).

S druhým problémom sa musí každý jazyk vysporiadať sám vytvorením vlastného „vzorkovníka“. Pre češtinu ho vytvoril Láda Lhotka a české delenie funguje celkom úspešne. I napriek tomu, že naše jazyky sú blízke, gramatické jemnosti ma viedli k pokusom s vytváraním slovenských vzorov pre delenie (a nie poslovenčením českých, čo by sa na prvý pohľad mohlo zdať jednoduchšie).

Rozdeľovanie slov na počítači sa nedá urobiť dokonale. Správnosť rozdelenia často závisí od kontextu, napr. na-drobiť a nad-robiť. Autorom rozdeľovacieho algoritmu pre \TeX , ako i „vzorkovníka“ pre angličtinu je Frank M. Liang. Svoje výsledky (hlavne problém vytváranie vzorov) zhrnul do Ph. D. Thesis (1983) na Standfordskej Univerzite. Navrhnutá metóda je dostatočne rýchla a láva možnosť pomocou iných vzorov učiť \TeX rozdeľovať i v iných jazykoch (napr. v slovenčine alebo v češtine). Kvalita delenia závisí potom už len od nás, koľko času venujeme príprave vzorov.

\TeX sa pri delení slova najskôr pozrie do slovníka výnimiek $\backslash\text{hyphenation}\{\dots\}$. Ak tam slovo nájde, rozdelí ho spôsobom, ako mu predpisujeme, napr. $\backslash\text{hyphenation}\{\text{výnim-ka}\}$. V takomto prípade už nepoužíva delenie pomocou „vzorkovníka“ (takže uvedením slova do bez rozdeľovníkov možno zabrániť rozdeľovaniu slova. v celom dokumente). Obdobne je tomu i u slov, v ktorých sa zadáva delenie explicitne pomocou $\backslash-$, napr. $\text{po}\backslash\text{-mlč}\backslash\text{-ka}$.

Pokiaľ nenastane ani jedna zo spomínaných možností, spustí Liangov algoritmus. Ukážme si, ako bude rozdelené slovo *Czechoslovakia* s anglickými vzormi zo súboru *enhyph.tex* (štandardne distribuovaného s $\zeta\text{\TeX}$ om). Pri rozdeľovaní slov veľké a malé písmen nehrajú rolu, preto \TeX najskôr „preloží“ slovo do malých písmen (samozrejme tento preklad musí byť presne definovaný). Na začiatok a koniec slova pridá extra znaky – bodky, ktoré budú v ďalšom signalizovať, že sa jedná o špeciálnu časť slova (začiatok a koniec).

² Každý nech si vyberie podľa svojho zmýšľania, či politickej príslušnosti.

Označované slovo

.czechoslovakia.

(dĺžky 16) rozdelí na všetky podslová dĺžky 1, 2, ... Takto dostáva nasledujúce podslová:

dĺžky 1: ., c, z, e, ...

dĺžky 2: .c, cz, ze, ...

dĺžky 3: .cz, zec, ech, ...

...

dĺžky 15: .czechoslovakia, czechoslovakia.

dĺžky 16: .czechoslovakia.

Potom pre každé obdržané podslovo nájde medzi vzormi jeho „ohodnotený“ (celými čísla od 0 do 9) ekvivalent (nemusi vždy existovať). V našom prípade nájde v spomínanom súbore nasledujúce „ohodnotené“ podslová:

cz₄, 2ze, 2ch, 3cho₂, os₄l, 4lov, 1va

podľa ktorých priradí hodnotu každému znaku v rozdeľovanom slove spôsobom:

- znaku, za ktorým nenasleduje žiadne číslo – priradí 0,
- znaku, za ktorým sa stretne viac čísel – priradí najväčšie z nich.

Chceme napríklad určiť hodnotu znaku e v spomínanom slove. Podľa vzorky 2ch tam prislúcha 2 podľa 3cho₂ hodnota 3, ale podľa cze₄ hodnota 4. Výsledná hodnota znaku e je teda 4. Uvedeným spôsobom určíme hodnoty všetkých znakov, čím dostaneme (0 nepíšeme)

cz₄e₄cho₂s₄lo₁v₁kia.

A *prekvapivý záver*: T_EX rozdelí slovo len za znakmi, ktoré majú nepárnu hodnotu. V našom prípade deliace miesto je teda jediné

czechoslo-vakia.

Pritom sa ešte kontrolujú momentálne nastavené hodnoty \lefthyphenmin (resp. \righthyphenmin) udávajúce kolkto prvých (resp. posledných) znakov nesmie byť v slove rozdelených: pre angličtinu býva \lefthyphenmin=2 a \righthyphenmin=3. (Toto sa ovšem nevzťahuje na slová rozdeľované pomocou \hyphenation{...}).

Presvedčiť sa, že T_EX skutočne rozdelí uvedené slovo daným spôsobom môžeme napísaním \showhyphens{Czechoslovakia} do svojho dokumentu a možné rozdelenia slova budú zapísané do súboru LOG. (Možno použiť i na celý paragraf.)

Viac informácií o niektorých jemnostiach pri rozdeľovaní slov v T_EXu môžu zvedavejší čitatelia nájsť v „biblii“ T_EXbook [1, str. 449–455. Nájdete tam tiež poznámku o tom, že vytváranie „vzorkovníkov“ je dobre platená práca pre expertov...

Napriek tomu, že medzi spomínaných expertov nepatrím, každodenné „T_EX-ovanie“ ma priviedlo k nutnosti vytvoriť aspoň nejakú slovenských vzorov pre delenie. Na rozdiel od Lianga, ktorý základ svojich vzorov vyrobil z obrovského množstva slov a ich odvolení (k dispozícii mal Webster’s Pocket Dictionary v elektronickej forme – zhruba 50 000 slov), použitá metóda je založená priamo na prepise gramatických pravidiel na rozdeľovanie slov. Možno niekomu poslúži ako inšpirácia na vytvorenie dokonalejšej verzie (prípadne poskytnite mi nápady na vylepšenie), alebo ako ukážka toho, kadiaľ cesta nevedie.

Prvá fáza spočíva v prepise mechanických pravidiel rozdeľovania (t.j. pravidiel typu – ak je medzi dvoma samohláskami jedna spoluhláska, ...). Tieto dávajú nie vždy uspokojivé výsledky pri príponách a predponách, čo však je len chybou krásy a

nie gramatickou! V [2, str. 53] sa totiž píše, že ak si rozhranie medzi predponou, slovným základom a príponou neuvedomujeme, pripúšťa sa i mechanické rozdeľovanie slov. (Teda obe rozdelenia „náj-dem“ i „ná-jdem“ sú správne.) V prípade mechanického rozdeľovania slov vznikajú problémy so zloženými slovami, napr. rozdelenie Československo je hrubou (gramatickou) chybou.

Pre slovenčinu obe hodnoty $\backslash\text{lefthyphenmin}$ a $\backslash\text{righthyphenmin}$ nastavíme na 2. Podľa [2] totiž pri rozdeľovaní neoddeľujeme obyčajne koncovú a začiatočnú slabiku, ak obsahuje iba samohlásku.

Na začiatku zabránime deleniu dvojhlások i_2e , i_2u , i_2a a spoluhlások d_2z , $d_2ž$, c_2h . Potom pridáme každú samohlásku samostatne s 1:

$a_1, i, \check{a}_1, e_1, \acute{e}_1, \dots$

Tým si zabezpečíme správne delenie „pekných“ slov, v ktorých sa pravidelne strieda spoluhláska so samohláskou (slov typu „ko-lá-č“ máme ošetrené s $\backslash\text{righthyphenmin}$). Ďalšie pravidlo: ak sú medzi samohláskami dve spoluhlásky (r, l počítam vždy medzi spoluhlásky) prepíšeme do tvaru:

$z b_1 b, z b_1 t, z b_1 \check{c}. \dots$

pre všetky možné dvojice spoluhlások, vrátane dz , $dž$, ch . *POZOR*: dvojice $z c_1 h$, $z d_1 z$, $z d_1 ž$ musíme zo zrejmych dôvodov vyčiarknuť zo „vzorkovníka“.

Pridaním ďalšieho pravidla – výskyt troch spoluhlások medzi dvoma samohláskami pridáme do slovníka „vhodné“ trojice v ohodnotení:

$\dots, s_3 t_2 r$ (ses-tra), $n_3 d_2 r$ (han-dra), $n_3 s_2 k$ (pán-sky), \dots

Problém výberu takýchto trojíc je zložitejší ako v predchádzajúcom prípade dvojíc, lebo nemôžeme pridať všetky možné trojice (a to nielen kvôli kvantite). Napr. pridaním trojice $t_3 \check{I} z k$ by sme pokazili správne rozdelenie slova otlkať, ktoré dostaneme cez dvojice. Problém súvisí s neslabičnými l , \acute{I} , r , $ž$, ktoré v týchto prípadoch plnia funkciu samohlásky. Problém sa teda sústreďuje na získanie zoznamu všetkých možných trojíc spoluhlások, vyskytujúcich sa v slovenských slovách. Nejaké štatistické výsledky možno nájsť v [3], rozhodne však nie sú uspokojivé.

Analogicky sa postupuje pri 4, 5, \dots spoluhláskach, vyskytujúcich sa spolu. Problém vyhľadávania a ohodnocovania takýchto „zhlukov“ je podobného charakteru ako v predchádzajúcom prípade.

$\dots, r_2 s_3 t_2 v, \dots$

V druhej fáze sa sústreďujeme na predpony. Doterajší prepis zabráňuje rozdeľovaniu niektorých predpôn, napr. v slove „vystáť“ (po predpone vy-). Slovo bude rozdelené vys-táť. Nasleduje teda pridávanie predpôn

$\dots, .d_0 z_3 k_4 r, .d_0 z_3 k_4 l, \dots$

Zdôrazňujem, že uvedenie častí slovného základu (niekedy i viac než jednopísmenovej), je nutné. Ináč by sme pripustili nesprávne delenia, napr. delenie slova do-ktor.

Súčasnou tejto fázy je i pridávanie najfrekventovanejších častí zložených slov (v podstate čosi ako slovná predpona):

$\dots, v_5 i a c_3 h_4$ (viac-hlasný), $s_5 \check{c} e s k o_5 s_4$, $.p_0 d_3 z$, $.š t v o r_3 r \dots$

i s uvedením častí slovného zkladu z rovnakého dôvodu ako v predšlom.

V tretej fáze sa zameráme na prípony. Ich „pridanie“ do úvah spočíva v úprave ohodnotenia „zhlukov“, prípadne ich presnejšej špeciifikácii v spojení s väčšou časťou slovného základu.

Nakoniec nasleduje pridávanie pravidiel pre cudzie slová

$e_1 a_2, u_1, 2, \dots$

a výnimiek do `\hyphenation{médi-um, ...}`. Tejto časti však nebola zatiaľ venovaná veľká pozornosť.

Čo dodať na záver: Snáď zodpovedať aktuálnu otázku „Ako bude rozdelené (slovo) Československo?“ Vytvoreným slovenským slovníkom: Čes-ko-slo-ven-sko ...

Literatura

- [1] Knuth D. E., The \TeX book, Addison-Wesley Pub. Comp. A AMS, 1986.
- [2] Oravec J. a Laca V., Príručka slovenského pravopisu pre školy, SPN, Bratislava, 1976.
- [3] Mistrík J., Frekvencia tvarov a konštrukcií v slovenčine, VEDA, Bratislava, 1985.

(Jana Chlebíková)

$\mathcal{A}\mathcal{M}\mathcal{S}$ - \LaTeX verze 1.0 a 1.1

V červenci roku 1990 nabídla American Mathematical Society veřejnosti první verzi $\mathcal{A}\mathcal{M}\mathcal{S}$ - \LaTeX u. Tato verze byla označena jako 1.0. V současné době je k dispozici verze 1.1, která se však od verze 1.0 liší zcela nepodstatně — ve verzi 1.1 byly pouze opraveny některé známé chyby. Dále budu proto jednoduše mluvit o $\mathcal{A}\mathcal{M}\mathcal{S}$ - \LaTeX u.

Pokud nepočítám příznivce PLAINu a odpůrce velkých balíků maker vůbec, stojí každý budoucí uživatel \TeX u před volbou, zda používat $\mathcal{A}\mathcal{M}\mathcal{S}$ - \TeX nebo \LaTeX (volba \LaTeX u je samozřejmě ta nejlepší!). Pravděpodobně ve snaze tuto volbu ulehčit, byly činěny různé pokusy, jak $\mathcal{A}\mathcal{M}\mathcal{S}$ - \TeX přiblížit \LaTeX u nebo naopak. Výsledkem nejznámějších dvou z těchto pokusů jsou balíky $\mathcal{L}\mathcal{A}\mathcal{M}\mathcal{S}$ - \TeX a $\mathcal{A}\mathcal{M}\mathcal{S}$ - \LaTeX . Můj soukromý dojem však je, že uživatelé \LaTeX u nebo $\mathcal{A}\mathcal{M}\mathcal{S}$ - \TeX u se pouze upevní ve víře, že jejich volba byla ta správná. Je třeba zdůraznit, že $\mathcal{L}\mathcal{A}\mathcal{M}\mathcal{S}$ - \TeX a $\mathcal{A}\mathcal{M}\mathcal{S}$ - \LaTeX není jedno a totéž a tyto systémy jsou dokonce nekompatibilní. Systém $\mathcal{L}\mathcal{A}\mathcal{M}\mathcal{S}$ - \TeX je rozšíření $\mathcal{A}\mathcal{M}\mathcal{S}$ - \TeX u, což mimo jiné vychází ze skutečnosti, že autorem $\mathcal{L}\mathcal{A}\mathcal{M}\mathcal{S}$ - \TeX u je M. D. Spivak (tj. autor $\mathcal{A}\mathcal{M}\mathcal{S}$ - \TeX u). $\mathcal{L}\mathcal{A}\mathcal{M}\mathcal{S}$ - \TeX přijímá mnohé rysy \LaTeX u jako např. automatické číslování formulí, kapitol, sekcí atd., možnost křížových referencí, okolí pro snadnou tvorbu tabulek atd. Věc, kterou $\mathcal{L}\mathcal{A}\mathcal{M}\mathcal{S}$ - \TeX od \LaTeX u nepřebírá, je okolí `picture` (okolí pro tvorbu jednoduchých obrázků přímo ve zdrojovém textu).

Na druhé straně dává $\mathcal{L}\mathcal{A}\mathcal{M}\mathcal{S}$ - \TeX velké možnosti při vytváření komutativních diagramů; kdo v textu potřebuje mnoho komutativních diagramů, tomu bych doporučil $\mathcal{L}\mathcal{A}\mathcal{M}\mathcal{S}$ - \TeX . Další výhodou $\mathcal{L}\mathcal{A}\mathcal{M}\mathcal{S}$ - \TeX u je poměrná uzavřenost systému a podrobná dokumentace (M. D. Spivak, který již dříve napsal *The Joy of \TeX* — manuál $\mathcal{A}\mathcal{M}\mathcal{S}$ - \TeX u — napsal podrobný manuál pro $\mathcal{L}\mathcal{A}\mathcal{M}\mathcal{S}$ - \TeX).

Na rozdíl od $\mathcal{L}\mathcal{A}\mathcal{M}\mathcal{S}$ - \TeX u je $\mathcal{A}\mathcal{M}\mathcal{S}$ - \LaTeX rozšířením \LaTeX u. Přitom rozšířením v tom smyslu, že základem je dokonce přímo soubor `latex.tex` bez nějakých úprav (předpokládá se, že bude použit soubor `latex.tex` verze 2.09 a to ne starší než z května 1986). Tato skutečnost mimo jiné znamená, že v $\mathcal{A}\mathcal{M}\mathcal{S}$ - \LaTeX u budou k dispozici všechny příkazy \LaTeX u tak, jak je uživatel zná z L^amportova manuálu (pouze příkazy pro volbu fontů budou mít poněkud jiný význam jak se zmíním dále). Kdo